# Realizing Fully-Integrated, Low-Power, Event-Based Pupil Tracking with Neuromorphic Hardware

Federico Paredes-Vallés[1]    Yoshitaka Miyatani[2]    Kirk Y. W. Scheper[1]

[1]Sony Advanced Visual Sensing AG, Schlieren, Switzerland
[2]Sony Semiconductor Solutions Corporation, Atsugi, Japan

## Abstract

*Eye tracking is fundamental to numerous applications, yet achieving robust, high-frequency tracking with ultra-low power consumption remains challenging for wearable platforms. While event-based vision sensors offer microsecond resolution and sparse data streams, they have lacked fully integrated, low-power processing solutions capable of real-time inference. In this work, we present the first battery-powered, wearable pupil-center-tracking system with complete on-device integration, combining event-based sensing and neuromorphic processing on the commercially available Speck2f system-on-chip with lightweight coordinate decoding on a low-power microcontroller. Our solution features a novel uncertainty-quantifying spiking neural network with gated temporal decoding, optimized for strict memory and bandwidth constraints, complemented by systematic deployment mechanisms that bridge the reality gap. We validate our system on a new multi-user dataset and demonstrate a wearable prototype with dual neuromorphic devices achieving robust binocular pupil tracking at 100 Hz with an average power consumption below 5 mW per eye. Our work demonstrates that end-to-end neuromorphic computing enables practical, always-on eye tracking for next-generation energy-efficient wearable systems.*

## 1. Introduction

Eye tracking is fundamental to a wide range of applications, from enhancing user interfaces and enabling medical diagnostics to advancing AR/VR, assistive technologies, and neuroscience [2, 18, 23, 38, 44, 49]. These applications demand precise, rapid estimation of eye position to ensure natural interaction and effective analysis of visual attention [3, 28, 32]. However, achieving robust, high-frequency tracking with low power consumption remains a significant challenge, particularly for wearable platforms [9, 42, 56].

Traditional eye-tracking systems typically rely on frame-based cameras and computationally intensive algorithms.
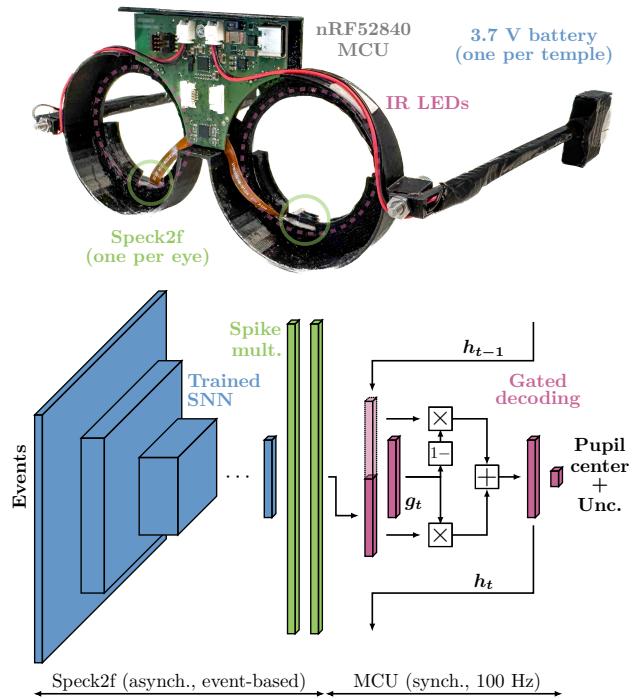


Figure 1. *Top*: Prototype of our fully-integrated, battery-powered pupil-center-tracking system based on neuromorphic technology, featuring two Speck2f devices, IR LEDs, and an nRF52840 MCU. *Bottom*: Our on-chip SNN processes input events asynchronously, while an MCU runs our off-chip decoding at 100 Hz to yield pupil-center and uncertainty estimates.

While capable of high spatial and moderate temporal resolution, these systems suffer from substantial power requirements and latency [9, 17, 45]. The continuous capture and processing of dense frames leads to inefficiencies unsuitable for lightweight, battery-powered devices. Furthermore, they are vulnerable to motion blur and limited dynamic range, which degrade performance during rapid eye movements or under challenging lighting conditions [4, 24, 27].

1

Event-based vision sensors (EVS) have recently emerged as a promising alternative. Unlike conventional image sensors, an EVS asynchronously detects per-pixel changes in log-brightness at microsecond resolution, producing sparse data streams that encode dynamic scene information [19]. This enables low-latency, high-frequency tracking and robustness to motion blur, while dramatically reducing power consumption [10, 11, 54]. However, the asynchronous and sparse nature of event data demands fundamentally different processing approaches, yet most existing methods rely on dense event representations and processing, limiting efficiency and scalability [10, 34, 40, 52, 54].

Achieving high-frequency, low-power eye tracking with EVS requires not only efficient sensing but also tightly integrated, event-driven processing. Neuromorphic computing offers a promising avenue through spiking neural networks (SNNs), which process information via discrete spikes, potentially enabling real-time computation at ultra-low power [13, 41, 53]. The sparse and asynchronous nature of both event data and spike-based processing makes neuromorphic systems naturally suited for EVS applications. However, existing solutions targeting deployment on dedicated neuromorphic processors rely on power-hungry development kits, offload computation, or fail to fully exploit hardware capabilities [8, 21, 22, 35–37].

In this work, we introduce the first fully integrated, battery-powered, wearable pupil-center-tracking system based on neuromorphic hardware. Leveraging the Speck2f system-on-chip (SoC) [41, 53], which integrates an EVS with dedicated SNN cores, we deploy a compact SNN that regresses pupil coordinates while respecting strict hardware constraints. By coupling event-based sensing, spike-based processing, and a low-power microcontroller (MCU) for efficient output decoding, we achieve robust, high-frequency tracking at very low power. Our contributions are:

- A novel uncertainty-quantifying SNN with gated temporal decoding and deployment pipeline optimized for the Speck2f's strict resource and readout constraints.
- The first fully integrated, battery-powered, wearable system for real-time, binocular pupil-center tracking combining event-based sensing and neuromorphic processing on dual Speck2f SoCs.
- Comprehensive validation on a new multi-user dataset and on a wearable prototype, demonstrating robust, high-frequency operation and low system-level power consumption.

## 2. Related work

Event-based eye tracking has gained significant attention due to the unique advantages of the EVS, such as low latency and high dynamic range. Early approaches include classical rule-based pipelines [4, 46], while recent work focuses on lightweight deep learning models [10, 15, 39, 50,

54, 55]. Notably, prior work demonstrated gaze tracking at kilohertz rates enabled by the microsecond resolution of the EVS [4, 46], and emphasized the importance of capturing temporal dynamics via advanced recurrent architectures [50]. Despite these advances, most solutions require substantial computational resources or offline processing, limiting their suitability for wearable applications.

To address these limitations, neuromorphic computing has emerged as a promising paradigm for event-based processing. Several works have explored SNNs for event-based eye tracking [20, 25], but despite achieving high accuracy, these solutions remain undeployed on actual neuromorphic hardware. The sole exception is Bonazzi *et al.* [6], who also tackled pupil-center tracking by deploying an SNN on the Speck2f SoC [41, 53]. However, their system relied on a power-hungry development kit and lacked support for continuous, real-time tracking due to network state resets between forward passes during training. Fully integrated, real-time, low-power event-based eye tracking on neuromorphic hardware therefore remains an open challenge.

Beyond eye tracking, other event-based solutions have been developed for neuromorphic hardware targeting applications like face detection [8], optical flow estimation [16, 37], and large language models [1], among many others [12, 22, 29, 31, 33, 43]. While these solutions demonstrate the potential for energy-efficient, real-time inference, most implementations rely on external computers for computation or interfacing, which undermines overall system-level power savings and limits practical deployment in battery-powered operation scenarios.

In this work, we address these limitations by achieving full system integration: all sensing and processing are performed on-device within a single, wearable platform that combines event-based sensing, sparse asynchronous processing, and a low-power MCU interface. Uniquely, our system deploys a novel SNN on two Speck2f SoCs, one per eye, enabling simultaneous, real-time pupil-center tracking with low system-level power consumption, hence demonstrating the practical viability of neuromorphic hardware for always-on wearable applications.

### 2.1. SynSense Speck2f

The core hardware in this work is the SynSense[1] Speck2f, a commercially available digital SoC that integrates a $128 \times 128$ EVS array with nine asynchronous convolutional SNN cores on a single die [41, 53]. This tight integration enables efficient, low-power, event-driven sensing and processing, making it well suited for always-on applications.

A key feature of the Speck2f is its fully asynchronous, event-driven processing. Unlike other digital neuromorphic architectures that use global synchronization to trigger neuron updates [13, 14], the Speck2f updates its spiking neu-
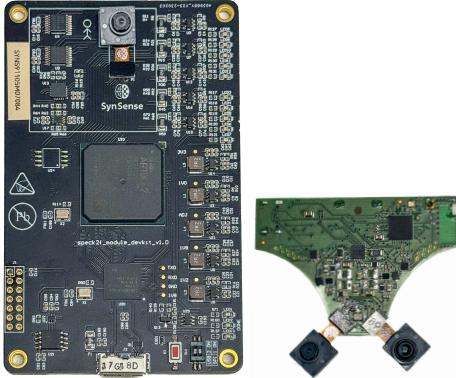
---

Figure 2. Speck2f support boards used in this work. *Left*: Speck2f development kit, featuring a Xilinx Artix-7 FPGA. *Right*: our custom PCB designed to accommodate two Speck2f units and featuring an nRF52840 MCU. The proportions of the boards are maintained to provide an accurate comparison of their sizes.

rons only when new input events arrive. This design allows each core to process spikes independently and immediately, routing internal spikes between cores via a network-on-chip without global coordination. Each core implements all neurons for a given convolutional layer, and processes input spikes in a first-in, first-out manner, performing (i) 2D convolution, (ii) spiking neuron update, and (iii) optional sum pooling for each spike.

Regarding the spiking model, the Speck2f supports leaky integrate-and-fire neurons [41, 47, 48, 53], where each neuron integrates incoming spikes with a linear leakage term and emits at most one output spike per update if the state exceeds a predefined threshold. In practice, leakage is applied synchronously via an external clock, which is at odds with the chip's otherwise asynchronous design. As a result, most prior work omits the leakage mechanism [7, 8, 41, 53], treating neurons as integrate-and-fire (IF) units without leakage. This makes neuron states dependent solely on the order of spikes, not their precise timing.

Each core operates under strict resource constraints, including 8-bit weights, 16-bit internal states, and limited bandwidth. The bandwidth of a core refers to the maximum number of synaptic operations[2] per second (SOPs) that can be processed without introducing latency. Most cores support up to 30M SOPs, while the first core, which is directly connected to the EVS, supports up to 100M SOPs [53].

For interfacing, SynSense mainly provides a development kit based on a Xilinx Artix-7 FPGA (see Fig. 2, left). All published work on this SoC to date has used this kit [6–8, 22, 41], which enables high-speed communication with a host computer via SynSense's proprietary protocol but incurs high static power consumption (approx. 0.6 W), far ex-

[2]A synaptic operation includes all computational steps required to process an incoming spike.

ceeding the Speck2f's own static power draw (approx. 0.5 mW), preventing low system-level power consumption.

The only alternative is a serial peripheral interface (SPI) that communicates directly with the Speck2f's dedicated readout core [41, 53]. This core, driven by an external slow clock (SCLK), supports up to 16 output neurons (with the first remaining silent), each with a dedicated processing engine for temporal averaging of spike counts over configurable windows. The readout core can report either the index and activity of the most active neuron or the activity of a specific neuron. Notably, most existing approaches employ more than 15 active output neurons [6–8, 22, 53], exceeding the hardware limits of the on-chip readout and rendering them incompatible with this native, low-power interface.

## 3. Method

The goal of this work is to develop the first fully integrated, low-power system for continuous, real-world pupil-center tracking using event-based sensing and processing. We design a specialized network architecture that fits the strict memory constraints of the Speck2f, along with optimized training and deployment pipelines. We demonstrate our solution by running it in parallel on two Speck2f devices connected to a custom circuit board with a low-power MCU, achieving efficient battery-powered operation integrated in a wearable prototype. The following sections detail each component of our framework.

### 3.1. Input format for off-chip training

For an ideal EVS [19], each event is generated at a pixel location and time when the log-brightness change since the last event at that pixel exceeds a polarity-specific threshold.

In Speck2f, events are sent directly from the sensing to the processing cores as they occur, without synchronization or accumulation [41, 53]. To avoid simulating microsecond-level event processing during training, we discretize event streams into non-overlapping time windows of duration $\Delta t$, as in [8, 21, 36, 37, 53]. Within each window, multiple events at the same pixel are aggregated into two-channel count images (one per polarity), providing a dense event representation. In the limit as $\Delta t \to 0$, this discretized approach approximates fully asynchronous event processing.

### 3.2. Network architecture and neuron model

A schematic of our model is shown in Fig. 1 (bottom). The core SNN employs $L = 7$ of the nine available Speck2f processing cores to form a feed-forward convolutional network comprising just 46.2k parameters. This network maps the two-channel $128 \times 128$ EVS input to 15 active output neurons, thereby matching the SPI-based readout limit. Each layer applies a strided $3 \times 3$ convolution, with channel counts per layer of {4, 12, 18, 27, 40, 60, 15}, chosen to fit core memory and feature map constraints [41, 53]. Similar

to the linear leakage (see Section 2.1), bias terms are omitted to eliminate the overhead associated with synchronous neuron updates [41, 47, 48, 53].

For the activation function, we use IF neurons optimized for compatibility with the Speck2f:

$$\tilde{v}_i[t] = v_i[t-1] - s_i^{\text{out}}[t-1]v_{\text{th}} + \sum_j w_{ij}s_{ij}^{\text{in}}[t] \quad (1)$$

$$v_i[t] = \max(v_{\text{min}}, \tilde{v}_i[t]) \quad (2)$$

$$s_i^{\text{out}}[t] = \begin{cases} 1, & \text{if } v_i[t] > v_{\text{th}} \text{ and on-chip} \\ \lfloor v_i[t]/v_{\text{th}} \rfloor, & \text{if } v_i[t] > v_{\text{th}} \text{ and off-chip} \\ 0, & \text{otherwise} \end{cases} \quad (3)$$

Here, $v_i[t]$ denotes the internal state, commonly referred to as *voltage*, of the $i$th neuron. This state evolves over time in response to weighted input spikes $s_i^{\text{in}}$ and the output spikes $s_i^{\text{out}}$ it generates, clipped below by $v_{\text{min}} = -10$ to mimic on-chip state quantization. On-chip, a neuron emits a single spike if its voltage exceeds the threshold $v_{\text{th}} = 1$. Off-chip, we model the neuron to emit multiple spikes per timestep proportional to $v_i[t]/v_{\text{th}}$, reducing the reality gap between training with discrete time windows $\Delta t$ and deployment on asynchronous event-based hardware [8, 53].

Additionally, we implement a soft reset [5]: after each output spike, the voltage is reduced by $s_i^{\text{out}}v_{\text{th}}$, as supported by Speck2f [41, 53]. To ensure correct behavior, we constrain synaptic weights to remain below $v_{\text{th}}$ via projected gradient descent. Otherwise, the state may remain above $v_{\text{th}}$ after an output spike, potentially increasing the reality gap. For training, we use surrogate gradients (periodic version of the inverse tangent derivative [8, 51]) due to the non-differentiable spike generation in Eq. 3, and, unlike [6], reset neuron states only after each sequence, thus supporting continuous operation.

### 3.3. Gated decoding: Regression from spikes

To derive continuous pupil-center coordinates from the 15 output spiking neurons, we propose a lightweight, stateful off-chip decoding (see Fig. 1, bottom) suitable for low-power MCUs. At each timestep, the output spiking activity $x_{L,t}$ is concatenated with a hidden state of equal length $h_{t-1}$. This vector passes through a fully connected layer (with learnable weights $W_g$ and biases $b_g$) with sigmoid activation $\sigma(\cdot)$ to compute per-neuron gating values $g_t$, which control the update of the hidden state:

$$g_t = \sigma(W_g[x_{L,t}, h_{t-1}] + b_g) \quad (4)$$

$$h_t = g_t \odot x_{L,t} + (1 - g_t) \odot h_{t-1} \quad (5)$$

The updated hidden state is min-max normalized and passed through a final linear layer to predict pupil-center coordinates (and aleatoric uncertainty, see Section 3.4):

$$h_t^{\text{norm}} = \frac{h_t - \min(h_t)}{\max(h_t) - \min(h_t)} \quad (6)$$

$$[\hat{x}, \hat{y}]^{\top} = \sigma(W_{xy}h_t^{\text{norm}} + b_{xy}) \quad (7)$$

$$\hat{u} = W_u h_t^{\text{norm}} + b_u \quad (8)$$

Unlike direct [8, 21, 37] or temporally-weighted spike decoding [6], our gating mechanism adaptively integrates temporal information, improving robustness to variable pupil speeds and occlusions while remaining efficient (approx. 1.1k floating-point operations). We avoid more complex recurrent units to minimize off-chip computation.

### 3.4. Training details

We train our network to regress pupil-center coordinates from event-camera data using a loss function that combines uncertainty-aware tracking and activity regularization:

$$\mathcal{L}(x_0, \theta) = \mathcal{L}_{\text{track.}}(x_0, \theta) + \beta\mathcal{L}_{\text{reg.}}(x_0, \theta) \quad (9)$$

where $x_0$ and $\theta$ denote the input to the SNN and the model parameters, respectively, while $\beta = 100$ is a scaling factor.

**Uncertainty-aware tracking:** We use the loss from Kendall et al. [26], which extends mean squared error of the predicted pupil-center location with a learned log variance $\hat{u} := \log\hat{\sigma}^2$ to model aleatoric uncertainty:

$$\mathcal{L}_{\text{track.}}(x_0, \theta) = \frac{1}{2}\exp(-\hat{u})\|y - \hat{y}\|^2 + \frac{1}{2}\hat{u} \quad (10)$$

where $y$ and $\hat{y}$ are the ground-truth and predicted pupil-center locations, respectively.

This uncertainty term serves as both a dynamic discount factor and a regularizer, encouraging reliable predictions and preventing unbounded uncertainty. It improves robustness to occlusions (e.g., during blinking) and low-motion scenarios (i.e., low signal-to-noise ratio due to limited input events), and removes the need for costly ground-truth observability labels in real-world data.

**Activity regularization:** To ensure real-time operation within the Speck2f's processing capacity, activity regularization is needed. Prior work targeting deployment on this SoC has employed either direct spike count penalization [8] or target-based optimization [6], but both have drawbacks: the former is sensitive to the balancing factor $\beta$ in Eq. 9 and ignores hardware constraints, while the latter can force unnecessary activity during periods of low input event rates. Instead, our approach leverages Speck2f-specific bandwidth constraints, defined in terms of SOPs (see Section 2.1), to match the capacity of the processing cores. We penalize internal activity only when these hardware-derived thresholds are exceeded, minimizing the reality gap and preserving on-chip accuracy. This allows
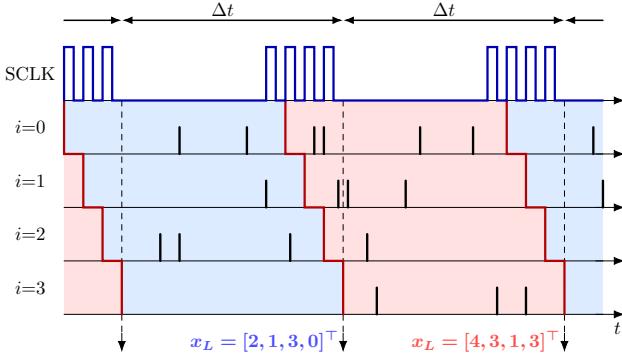
4

Figure 3. Cyclic readout strategy for a four-neuron layer. The MCU sequentially samples each neuron via SPI, synchronized to the SCLK signal, as rapidly as possible. Vertical black lines indicate spike events; colored regions show sampling windows. Dashed arrows mark the completion of each readout cycle, triggering the off-chip decoding process.
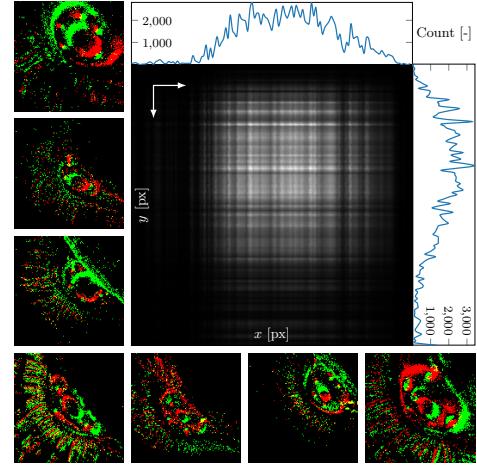


Figure 4. Representative samples from our real-world dataset (*left* and *bottom*) show event streams from different users and acquisition conditions. Colors used to distinguish event polarities: red for negative, green for positive. The *central* heatmap displays the distribution of ground-truth pupil centers, with marginal histograms indicating coverage along each axis.

the model to converge to the best possible tracking solution that respects these constraints.

For each layer $l$, SOPs are computed from the input activity and connectivity within a time window $\Delta t$:

$$\text{SOPs}_l = \frac{\boldsymbol{x}_{l-1} \, \boldsymbol{m}_{l-1,l}}{\Delta t} \qquad (11)$$

where $\boldsymbol{m}_{l-1,l}$ is the binary connectivity matrix between layers $l-1$ and $l$.

The activity regularization term is then defined as:

$$\mathcal{L}_{\text{reg.}}(\boldsymbol{x}_0, \theta) = \sum_{l=2}^{L-1} \kappa(\text{SOPs}_l, \text{SOPs}_l^{\text{th}}) + \kappa(\boldsymbol{x}_{\text{L}}, \boldsymbol{x}_{\text{L}}^{\text{th}}) \qquad (12)$$

$$\kappa(a, b) = \frac{\max(0, a - b)}{b} \qquad (13)$$

where the superscript "th" denotes the threshold above which excess activity is penalized. We use $\text{SOPs}_l^{\text{th}} = 20\text{M}$ and $\boldsymbol{x}_{\text{L}}^{\text{th}} = 83.3\text{k}$. We do not regularize the first layer, as its activity depends solely on sensor settings and scene dynamics, nor the off-chip gated decoding. The final layer's output spikes $\boldsymbol{x}_{\text{L}}$ are directly regularized to control readout load.

**Optimizer and other details:** During training, event-camera sequences are processed sequentially in batches of 32, with a discretization window of $\Delta t = 10$ ms. The loss is computed at each simulated timestep, and backpropagation is used to update the model parameters. We use the AdamW optimizer [30] with a weight decay of 0.05, as narrower weight distributions were empirically found to reduce the reality gap when deploying on the Speck2f SoC. Quantization-aware training is omitted as it provided negligible benefit for on-chip inference; instead, post-training quantization to 8-bit weights and 16-bit activations is applied before deployment.

### 3.5. The road to an integrated solution

Beyond our SNN and training procedure, additional mechanisms are required for a fully integrated solution where a low-power MCU interfaces two Speck2f devices via SPI. For more implementation-level details, including SPI commands structures, configuration encoding, and initialization steps, please refer to the supplementary material.

**Cyclic readout:** Retrieving the spiking activity from all output neurons over a time interval $\Delta t$ is a crucial preliminary step for our off-chip decoding process. However, as described in Section 2.1, a key challenge is the Speck2f readout core's limitation to 16 accessible neurons, which must be polled sequentially via SPI. We address this with a cyclic readout strategy (see Fig. 3): the MCU time-multiplexes the sampling of each neuron's activity within every $\Delta t$ window, synchronized to the SCLK signal, as rapidly as possible. This introduces minimal temporal offsets between readouts, mainly determined by the SCLK cycle duration (approx. 175 $\mu$s in our case). By cycling through all neuron indices, we construct output vectors $\boldsymbol{x}_{\text{L}}$ representing the final layer's spiking activity for off-chip decoding.

**Spike multiplier:** A notable limitation of the Speck2f's readout core is that, as in Section 2.1, it provides quantized moving averages over 16 SCLK cycles[3], rather than total spike counts. To bypass this moving average, simply increasing spike rates by scaling weights and thresholds is infeasible, as weights must remain below thresholds (see Section 3.2). To overcome this, we introduce an intermediate

---

[3]Available moving-average windows are 1, 16, or 32 SCLK cycles [47]. We use 16 to retrieve output spikes over $\Delta t$ in a single SPI transaction.
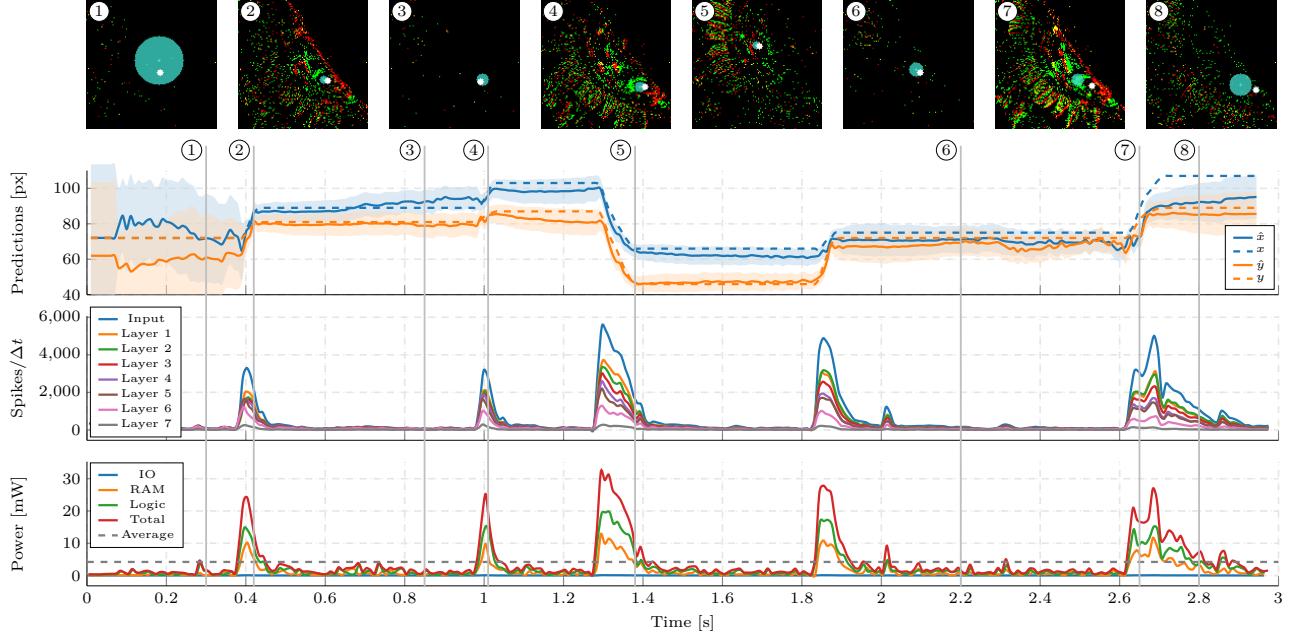
Figure 5. Example sequence from our real-world dataset and on-chip model evaluation. *Top*: Event frames at key time points (white: ground truth; blue: prediction). *Middle*: Predicted pupil coordinates with uncertainty, event counts, and spikes per layer. *Bottom*: Processing power consumption for each Speck2f subsystem, with a total average of 4.22 mW. Vertical lines indicate correspondence between frames and time-series data. Results obtained with the Speck2f development kit.

layer connecting each output neuron to $N$ neurons with unit weights and thresholds, effectively multiplying the spike count by $N$ at the cost of an additional core. We use $N = 4$ instead of 16 due to a known hardware issue in mapping neurons to the readout core [47], which necessitates adding a secondary layer with 4:1 connectivity (see supplementary material), thus effectively multiplying output activity by 16.

### 3.6. Wearable hardware overview

To realize a fully integrated pupil-center-tracking solution, we developed a custom printed circuit board (PCB) integrating two Speck2f devices (each with a 1.98 mm lens) and an nRF52840 MCU[4], interfaced via SPI at 100 Hz. The MCU manages power, device programming, real-time Bluetooth monitoring, and output decoding. The system supports USB or dual 3.7 V lithium coin battery power, enabling untethered wearable use. For robust detection under varying lighting, the PCB includes two infrared (IR) LED rings with six LEDs each[5]. Dedicated analog circuits monitor real-time power consumption of both Speck2f devices and the MCU. See Fig. 2 (right) for the custom PCB and Fig. 1 (top) for the wearable prototype.

---

[4]MCU selection optimization for lower system-level power consumption was beyond the scope of this work.

[5]Illumination optimization was beyond the scope of this work.

## 4. Experiments

Our primary objective is to achieve continuous pupil-center tracking with a fully integrated, low-power neuromorphic system. We do not directly compare to existing event-based methods, as prior work is generally unconstrained by hardware limitations and does not face the strict memory and bandwidth constraints of our system. Instead, we collected a dedicated real-world, multi-user dataset with the Speck2f for both training and validation, enabling rigorous evaluation under realistic operating conditions, which is then followed by an assessment using our wearable prototype.

### 4.1. Real-world dataset

Our dataset[6] contains 432 sequences from 8 users: 369 for training and 63 for validation. Each 3-second sequence captures data from a single eye as participants tracked a moving dot on a display at a fixed distance using a chinrest. To increase pupil size diversity, three background colors were used, and data was collected with four Speck2f devices to account for hardware variability. Natural eye behaviors, including blinks, were retained. Ground-truth pupil centers were manually annotated at 100 Hz. Fig. 4 shows representative samples and statistics. As shown, the Speck2f was angled to maximize eye region coverage. During training, data augmentation included random affine transformations.

---
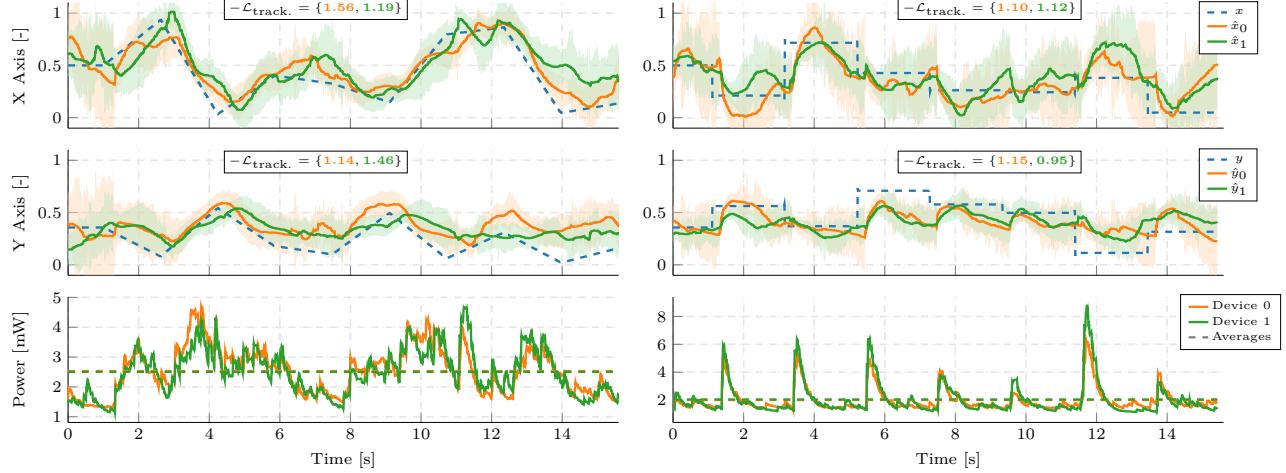
[6]To be made publicly available upon publication.

Figure 6. Wearable prototype results as a user performs smooth pursuit (left) and saccadic eye movements (right) while tracking a moving dot on a display at a fixed distance. *Top and middle rows*: Ground-truth and predicted pupil-center coordinates with uncertainty estimates for left (device 0) and right (1) eyes. *Bottom row*: Instantaneous power consumption per device and sequence averages. In both sequences, the MCU maintains an average power consumption of approx. 13.6 mW. Additional results can be found in the supplementary material.

## 4.2. Results on real-world dataset

We first evaluate our approach on the validation set of the real-world dataset, focusing on tracking accuracy and the reality gap. These experiments utilize the Speck2f development kit (see Fig. 2, left), which enables controlled event stream input and the monitoring of the internal activity of each layer. Quantitative results are presented in Table 1, with qualitative trends illustrated in Figs. 5 and 7.

Our gated decoding configurations achieve the best on-chip performance. Models predicting shared and per-axis uncertainty yield average L2 errors (lower is better, ↓) of 9.99 and 9.91 pixels, respectively, closely matching their GPU results with minimal reality gap. Given their comparable accuracy, we select the shared uncertainty variant to reduce off-chip computation. Notably, models with uncertainty estimation achieve up to 24% error reduction over no-uncertainty baselines. The performance advantage of gated over direct decoding (i.e., nearly 50% error reduction across uncertainty configurations) underscores the importance of temporal reasoning for this task, suggesting that neuromorphic hardware supporting more sophisticated recurrent dynamics could enable further accuracy improvements.

The selected model produces well-calibrated confidence estimates, evidenced by the correlation between predicted uncertainty and actual error in Fig. 7. On-chip median error increases from 3.74 pixels for the 10% most confident predictions to 7.27 pixels across all samples, showing that the predicted uncertainty effectively identifies challenging cases. See the supplementary material for additional results.

Finally, Fig. 5 highlights the energy efficiency and tracking performance of our solution over time. As shown, our predictions closely follow the ground-truth locations,

| Decoding | Uncertainty | L2 Error [px] ↓ | | SOPs [M/s] ↓ |
| --- | --- | --- | --- | --- |
| | | GPU | Speck2f | |
| Direct | No | 21.78 | 21.67 ± 0.04 | 8.73 |
| | Shared | 19.51 | 19.46 ± 0.02 | 15.33 |
| | Per-axis | 19.24 | 19.08 ± 0.01 | 15.01 |
| Gated | No | 13.11 | 13.24 ± 0.01 | <u>5.61</u> |
| | Shared | **9.74** | <u>9.99 ± 0.04</u> | **5.57** |
| | Per-axis | <u>9.80</u> | **9.91 ± 0.02** | 6.04 |

Table 1. Quantitative evaluation on our real-world dataset. Best in bold, runner up underlined. On-chip results are averaged over 5 runs to account for hardware stochasticity [41, 53]. Results obtained with the Speck2f development kit.
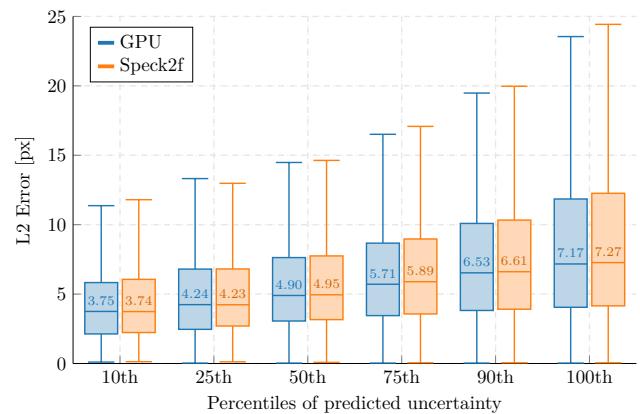


Figure 7. L2 error of the SNN model with gated decoding across percentiles of the shared predicted uncertainty for GPU-based and Speck2f-based (with development kit) evaluations.
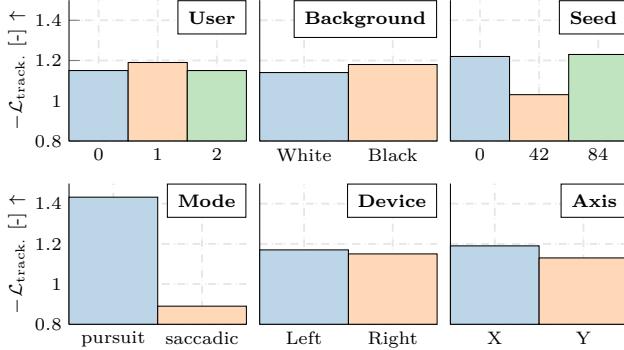
7

Figure 8. Tracking performance variability across conditions using our wearable prototype. Each subplot shows results for a specific condition, averaged over all other variables.

demonstrating high on-chip tracking accuracy. Furthermore, the Speck2f maintains consistently low power consumption, averaging 4.22 mW, which strongly correlates with both the input event rate and internal spike activity.

### 4.3. Results with wearable prototype

We demonstrate and evaluate our solution on our wearable prototype (see Fig. 1, top), focusing on tracking accuracy and real-time, low-power operation. Three users wore the prototype and tracked a moving stimulus on a display at a fixed distance using a chinrest. To align the predicted and ground-truth pupil locations, we recorded a calibration sequence for each user, fitting an affine transformation via least squares. We evaluate performance across six factors: user identity, background color, stimulus seed (controlling trajectory variation), tracking mode, device (i.e., left/right eye), and coordinate axis, reporting accuracy via the negative mean tracking error $-\mathcal{L}_{\text{track.}}$ (higher is better, ↑). Quantitative results are shown in Fig. 8, with qualitative trends illustrated in Fig. 6.

Our solution maintains consistent tracking performance across most conditions, demonstrating robust operation across users, backgrounds, devices, and axes (with slightly better performance along the x-axis), as shown in Fig. 8. The greatest variability appears in the stimulus seed and tracking mode. Different trajectory seeds produce varying difficulty levels based on speed, direction, and transition combinations. Additionally, saccadic movements yield considerably lower accuracy than smooth pursuit due to fundamental differences in eye movement dynamics and their interaction with event-based sensing. Smooth pursuit involves continuous, predictable tracking, generating sustained, high signal-to-noise ratio event streams. In contrast, saccades are rapid, ballistic eye movements producing bursty event patterns during transitions followed by minimal activity during fixations. This temporal sparsity challenges our architecture's limited temporal modeling capacity.

Fig. 6 illustrates these differences. During smooth pursuit (left), both eyes track closely with predicted coordinates following ground truth consistently and tight uncertainty estimates reflecting model confidence. Power consumption remains stable at around 2.5 mW per device. During saccadic movements (right), predicted trajectories show larger overall deviations from ground truth, with both uncertainty and errors increasing during the challenging fixation phases. Power consumption exhibits brief spikes of 6–8 mW during saccades, corresponding to event bursts from rapid eye movements, before returning to a baseline during fixations of approx. 1.5 mW. Together with Fig. 5, these power profiles validate the event-driven nature of our system: power scales with visual dynamics, demonstrating efficient energy usage that responds directly to eye movement complexity.

## 5. Limitations

While our approach demonstrates fully integrated, low-power, event-based pupil tracking, several opportunities for improvement remain. First, our model employs simple IF neurons and a basic CNN structure due to Speck2f's strict constraints. More sophisticated architectures could enhance temporal reasoning but would require more flexible neuromorphic hardware. Similarly, the proposed spike multiplier mechanism consumes two of the nine processing cores to bypass readout limitations, reducing resources for network complexity. Future neuromorphic hardware with direct spike count retrieval or a dedicated core for more flexible on-chip decoding would eliminate this overhead.

Second, our current 100 Hz operating frequency matches our ground-truth annotation rate. This represents an under-utilized system capability: the sparse, asynchronous nature of the SNN enables scaling to higher frequencies with negligible power increase, as processing remains event-driven and proportional to visual dynamics rather than output sampling rate. Given access to high-frequency ground truth for decoder training, our system could achieve significantly higher tracking rates with minimal additional power cost; a key scalability advantage of our neuromorphic approach.

Finally, the reported tracking errors stem from hardware constraints and dataset complexity, including multi-user variability, blinks, and rapid movements followed by long fixations. Despite these limitations, our solution suits applications requiring approximate gaze estimation, such as vergence estimation or coarse gaze-based interaction.

## 6. Conclusion

This work bridges the gap between neuromorphic hardware capabilities and practical deployment requirements for continuous wearable eye tracking. By addressing the unique challenges of tight integration, we demonstrate that this technology can power practical, always-on event-based vi-

sion systems with minimal energy overhead. Our SNN, powered by the gated decoding and uncertainty predictions, enables robust, continuous pupil-center tracking while respecting the stringent constraints of on-chip processing. The resulting fully integrated system, featuring two Speck2f SoCs, achieves sub-5 mW average power consumption per eye without sacrificing accuracy, demonstrating that neuromorphic processors can deliver practical performance beyond laboratory prototypes. Beyond eye tracking, the principles and methods introduced here provide a blueprint for intelligent, energy-efficient systems across diverse event-based computer vision tasks requiring always-on operation.

## Acknowledgments

## References

[1] Steven Abreu, Sumit Bam Shrestha, Rui-Jie Zhu, and Jason Eshraghian. Neuromorphic principles for efficient large language models on Intel Loihi 2. *arXiv:2503.18002*, 2025. 2

[2] Isayas Berhe Adhanom, Paul MacNeilage, and Eelke Folmer. Eye tracking in virtual reality: A broad review of applications and challenges. *Virtual Real.*, 27(2):1481–1505, 2023. 1

[3] Richard Andersson, Marcus Nyström, and Kenneth Holmqvist. Sampling frequency and eye-tracking measures. *J. Eye Movement Res.*, 3(3), 2010. 1

[4] Anastasios N. Angelopoulos, Julien N. P. Martel, Amit P. S. Kohli, Jorg Conradt, and Gordon Wetzstein. Event based, near eye gaze tracking beyond 10,000 Hz. *arXiv:2004.03577*, 2020. 1, 2

[5] Guillaume Bellec, Darjan Salaj, Anand Subramoney, Robert Legenstein, and Wolfgang Maass. Long short-term memory and learning-to-learn in networks of spiking neurons. *Adv. Neural Inform. Process. Syst.*, 31, 2018. 4

[6] Pietro Bonazzi, Sizhen Bian, Giovanni Lippolis, Yawei Li, Sadique Sheik, and Michele Magno. Retina: Low-power eye tracking with event camera and spiking hardware. In *IEEE Conf. Comput. Vis. Pattern Recog. Worksh.*, pages 5684–5692, 2024. 2, 3, 4

[7] Francesco Branca. Optical flow determination with integrate and fire neurons. MSc Thesis, Delft University of Technology, 2024. 3

[8] Caterina Caccavella, Federico Paredes-Vallés, Marco Cannici, and Lyes Khacef. Low-power event-based face detection with asynchronous neuromorphic hardware. In *Int. Joint Conf. Neural Netw.* IEEE, 2024. 2, 3, 4

[9] Marco Carminati, Filippo Melloni, Giulio Marano, Alberto Pettenella, Daniele Bani, Daniele Maria Crafa, Andrea Aspesi, Andrew Duchowski, Tommaso Ongarello, and Luca Merigo. Energy-aware benchmarking of wearable eye trackers. In *Eye Tracking Res. Appl.*, 2025. 1

[10] Qinyu Chen, Zuowen Wang, Shih-Chii Liu, and Chang Gao. 3ET: Efficient event-based eye tracking using a change-based ConvLSTM network. In *IEEE Biomed. Circuits Syst. Conf.* IEEE, 2023. 2

[11] Qinyu Chen, Chang Gao, Min Liu, Daniele Perrone, Yan Ru Pei, Zuowen Wang, Zhuo Zou, Shihang Tan, Tao Han, Guorui Lu, et al. Event-based eye tracking. 2025 Event-based vision workshop challenge survey. *arXiv:2504.18249*, 2025. 2

[12] Giulia D'Angelo, Victoria Clerico, Chiara Bartolozzi, Matej Hoffmann, P. Michael Furlong, and Alexander Hadjiivanov. Wandering around: A bioinspired approach to visual attention through object motion sensitivity. *arXiv:2502.06747*, 2025. 2

[13] Mike Davies, Narayan Srinivasa, Tsung-Han Lin, Gautham Chinya, Yongqiang Cao, Sri Harsha Choday, Georgios Dimou, Prasad Joshi, Nabil Imam, Shweta Jain, et al. Loihi: A neuromorphic manycore processor with on-chip learning. *IEEE Micro*, 38(1):82–99, 2018. 2

[14] Mike Davies, Andreas Wild, Garrick Orchard, Yulia Sandamirskaya, Gabriel A. Fonseca Guerra, Prasad Joshi, Philipp Plank, and Sumedh R. Risbud. Advancing neuromorphic computing with Loihi: A survey of results and outlook. *IEEE*, 109(5):911–934, 2021. 2

[15] Junyuan Ding, Ziteng Wang, Chang Gao, Min Liu, and Qinyu Chen. FACET: Fast and accurate event-based eye tracking using ellipse modeling for extended reality. *arXiv:2409.15584*, 2024. 2

[16] Julien Dupeyroux, Jesse J. Hagenaars, Federico Paredes-Vallés, and Guido C. H. E. de Croon. Neuromorphic control for optic-flow-based landing of mavs using the loihi processor. In *IEEE Int. Conf. Robot. Autom.*, pages 96–102. IEEE, 2021. 2

[17] Yu Feng, Tianrui Ma, Yuhao Zhu, and Xuan Zhang. Bliss-Cam: Boosting eye tracking efficiency with learned in-sensor sparse sampling. *arXiv:2404.15733*, 2024. 1

[18] Anke Fischer-Janzen, Thomas M. Wendt, and Kristof Van Laerhoven. A scoping review of gaze and eye tracking-based control methods for assistive robotic arms. *Front. Robot. AI*, 11, 2024. 1

[19] Guillermo Gallego, Tobi Delbrück, Garrick Orchard, Chiara Bartolozzi, Brian Taba, Andrea Censi, Stefan Leutenegger, Andrew J. Davison, Joerg Conradt, Kostas Daniilidis, and Davide Scaramuzza. Event-based vision: A survey. *IEEE Trans. Pattern Anal. Mach. Intell.*, 43(1):46–68, 2019. 2, 3

[20] Stijn Groenen, Marzieh H. Varposhti, and Mahyar Shahsavari. GazeSCRNN: Event-based near-eye gaze tracking using a spiking neural network. *arXiv:2503.16012*, 2025. 2

[21] Jesse J. Hagenaars, Federico Paredes-Vallés, and Guido C. H. E. De Croon. Self-supervised learning of event-based optical flow with spiking neural networks. *Adv. Neural Inform. Process. Syst.*, 34:7167–7179, 2021. 2, 3, 4

[22] Adam D. Hines, Michael Milford, and Tobias Fischer. A compact neuromorphic system for ultra–energy-efficient, on-device robot localization. *Sci. Robot.*, 10(103):eads3968, 2025. 2, 3

[23] Kenneth Holmqvist, Marcus Nyström, Richard Andersson, Richard Dewhurst, Halszka Jarodzka, and Joost van de Weijer. *Eye tracking: A comprehensive guide to methods and measures*. Oxford University Press, Oxford, 2011. 1

[24] Ignace TC Hooge, Diederick C Niehorster, Roy S Hessels, Jeroen S Benjamins, and Marcus Nyström. How robust are wearable eye trackers to slow and fast head and body movements? *Behav. Res. Methods*, 55(8):4128–4142, 2023. 1

[25] Yizhou Jiang, Wenwei Wang, Lei Yu, and Chu He. Eye tracking based on event camera and spiking neural network. *Electronics*, 13(14), 2024. 2

[26] Alex Kendall and Yarin Gal. What uncertainties do we need in bayesian deep learning for computer vision? *Adv. Neural Inform. Process. Syst.*, 30, 2017. 4, 12

[27] Rakshit Kothari, Zhizhuo Yang, Christopher Kanan, Reynold Bailey, Jeff B. Pelz, and Gabriel J. Diaz. Gaze-in-wild: A dataset for studying eye and head coordination in everyday activities. *Sci. Rep.*, 10(1), 2020. 1

[28] Alexander Leube, Katharina Rifai, and Siegfried Wahl. Sampling rate influences saccade detection in mobile eye tracking. *J. Eye Movement Res.*, 10(3), 2017. 1

[29] Qian Liu, Ole Richter, Carsten Nielsen, Sadique Sheik, Giacomo Indiveri, and Ning Qiao. Live demonstration: Face recognition on an ultra-low power event-driven convolutional neural network asic. In *IEEE Conf. Comput. Vis. Pattern Recog. Worksh.*, 2019. 2

[30] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *Int. Conf. Learn. Represent.*, 2019. 5

[31] Zifan Ning, Chaojin Chen, Xiang Cheng, Wangzi Yao, Tielin Zhang, and Bo Xu. Neuromorphic incremental on-chip learning with hebbian weight consolidation. *arXiv:2311.18340*, 2023. 2

[32] Marcus Nyström, Diederick C. Niehorster, Richard Andersson, and Ignace T. C. Hooge. The Tobii Pro Spectrum: A useful tool for studying microsaccades? *Behav. Res. Methods*, 53:335–353, 2021. 1

[33] Garrick Orchard, E. Paxon Frady, Daniel Ben Dayan Rubin, Sophia Sanborn, Sumit Bam Shrestha, Friedrich T. Sommer, and Mike Davies. Efficient neuromorphic signal processing with Loihi 2. In *IEEE Worksh. Signal Proc. Syst.*, pages 254–259. IEEE, 2021. 2

[34] Federico Paredes-Vallés and Guido CHE De Croon. Back to event basics: Self-supervised learning of image reconstruction for event cameras via photometric constancy. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 3446–3455, 2021. 2

[35] Federico Paredes-Vallés, Kirk Y. W. Scheper, and Guido C. H. E. de Croon. Unsupervised learning of a hierarchical spiking neural network for optical flow estimation: From events to global motion perception. *IEEE Trans. Pattern Anal. Mach. Intell.*, 42(8):2051–2064, 2019. 2

[36] Federico Paredes-Vallés, Kirk Y. W. Scheper, Christophe De Wagter, and Guido C. H. E. De Croon. Taming contrast maximization for learning sequential, low-latency, event-based optical flow. In *IEEE/CVF Int. Conf. Comput. Vis.*, pages 9695–9705, 2023. 3

[37] Federico Paredes-Vallés, Jesse J. Hagenaars, Julien Dupeyroux, Stein Stroobants, Yingfu Xu, and Guido C. H. E. de Croon. Fully neuromorphic vision and control for autonomous drone flight. *Sci. Robot.*, 9(90):eadi0591, 2024. 2, 3, 4

[38] Joseph R. Pauszek. An introduction to eye tracking in human factors healthcare research and medical device testing. *Hum. Factors Healthc.*, 3, 2023. 1

[39] Yan Ru Pei, Sasskia Brüers, Sébastien Crouzet, Douglas McLelland, and Olivier Coenen. A lightweight spatiotemporal network for online eye tracking with event camera. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 5780–5788, 2024. 2

[40] Henri Rebecq, René Ranftl, Vladlen Koltun, and Davide Scaramuzza. Events-to-video: Bringing modern computer vision to event cameras. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 3857–3866, 2019. 2

[41] Ole Richter, Yannan Xing, Michele De Marchi, Carsten Nielsen, Merkourios Katsimpris, Roberto Cattaneo, Yudi Ren, Yalun Hu, Qian Liu, Sadique Sheik, et al. Speck: A smart event-based vision sensor with a low latency 327k neuron convolutional neuronal network processing pipeline. *arXiv:2304.06793*, 2023. 2, 3, 4, 7

[42] Thiago Santini, Wolfgang Fuhl, and Enkelejda Kasneci. Pure: Robust pupil detection for real-time pervasive eye tracking. *Comput. Vis. Image Underst.*, 170:40–50, 2018. 1

[43] Sumit Bam Shrestha, Jonathan Timcheck, Paxon Frady, Leobardo Campos-Macias, and Mike Davies. Efficient video and audio processing with Loihi 2. In *IEEE Int. Conf. Acoust., Speech, Signal Process.*, pages 13481–13485. IEEE, 2024. 2

[44] Mohammed Tahri Sqalli, Begali Aslonov, Mukhammadjon Gafurov, Nurmukhammad Mukhammadiev, and Yahya Sqalli Houssaini. Eye tracking technology in medical practice: A perspective on its diverse applications. *Front. Med. Technol.*, 5, 2023. 1

[45] Niklas Stein, Diederick C. Niehorster, Tamara Watson, Frank Steinicke, Katharina Rifai, Siegfried Wahl, and Markus Lappe. A comparison of eye tracking latencies among several commercial head-mounted displays. *i-Perception*, 12(1), 2021. 1

[46] Timo Stoffregen, Hossein Daraei, Clare Robinson, and Alexander Fix. Event-based kilohertz eye tracking using coded differential lighting. In *IEEE/CVF Winter Conf. Appl. Comput. Vis.*, pages 2515–2523, 2022. 2

[47] Synsense. *Samna API documentation*, . 3, 4, 5, 6, 13, 15

[48] Synsense. *Speck development kit datasheet*, . Aug. 2023. 3, 4, 15

[49] Niilo V. Valtakari, Ignace T. C. Hooge, Charlotte Viktorsson, Pär Nyström, Terje Falck-Ytter, and Roy S. Hessels. Eye tracking in human interaction: Possibilities and limitations. *Behav. Res. Methods*, 53(4):1592–1608, 2021. 1

[50] Zhong Wang, Zengyu Wan, Han Han, Bohao Liao, Yuliang Wu, Wei Zhai, Yang Cao, and Zheng-Jun Zha. MambaPupil: Bidirectional selective recurrent model for event-based eye tracking. In *IEEE Conf. Comput. Vis. Pattern Recog. Worksh.*, pages 5762–5770, 2024. 2

[51] Philipp Weidel and Sadique Sheik. Wavesense: Efficient temporal convolutions with spiking neural networks for keyword spotting. *arXiv:2111.01456*, 2021. 4

[52] Yilun Wu, Federico Paredes-Vallés, and Guido CHE De Croon. Lightweight event-based optical flow estimation via iterative deblurring. In *IEEE Int. Conf. Robot. Autom.*, pages 14708–14715. IEEE, 2024. 2

[53] Man Yao, Ole Richter, Guangshe Zhao, Ning Qiao, Yannan Xing, Dingheng Wang, Tianxiang Hu, Wei Fang, Tugba Demirci, Michele De Marchi, et al. Spike-based dynamic computing with asynchronous sensing-computing neuromorphic chip. *Nat. Commun.*, 15(1):4464, 2024. 2, 3, 4, 7, 13

[54] Baoheng Zhang, Yizhao Gao, Jingyuan Li, and Hayden Kwok-Hay So. See: Co-designing a sub-millisecond latency event-based eye tracking system with submanifold sparse convolution. In *IEEE Conf. Comput. Vis. Pattern Recog. Worksh.*, 2024. 2

[55] Guangrong Zhao, Yurun Yang, Jingwei Liu, Ning Chen, Yiran Shen, Hongkai Wen, and Guohao Lan. Ev-eye: Rethinking high-frequency eye tracking through the lenses of event cameras. *Adv. Neural Inform. Process. Syst.*, 36:62169–62182, 2023. 2

[56] Lisen Zhu, Jianan Chen, Huixin Yang, Xinkai Zhou, Qihang Gao, Rui Loureiro, Shuo Gao, and Hubin Zhao. Wearable near-eye tracking technologies for health: A review. *Bioengineering*, 11(7):738, 2024. 1

# Realizing Fully-Integrated, Low-Power, Event-Based Pupil Tracking with Neuromorphic Hardware
## — Supplementary Material —

Federico Paredes-Vallés[1]     Yoshitaka Miyatani[2]     Kirk Y. W. Scheper[1]

[1]Sony Advanced Visual Sensing AG, Schlieren, Switzerland
[2]Sony Semiconductor Solutions Corporation, Atsugi, Japan

## S1. Uncertainty calibration

To validate our uncertainty estimates [26], we analyze its calibration using reliability diagrams adapted for regression. For each prediction $\hat{y}$ of pupil center location with estimated standard deviations $\hat{\sigma} = (\hat{\sigma}_x, \hat{\sigma}_y)$, we compute the normalized error via the Mahalanobis distance:

$$z = \sqrt{(y - \hat{y})^\top \Sigma^{-1}(y - \hat{y})}, \qquad (14)$$

where $\Sigma = \mathrm{diag}(\hat{\sigma}_x^2, \hat{\sigma}_y^2)$ is the diagonal covariance matrix.

Under a bivariate Gaussian assumption with diagonal covariance, $z^2$ follows a $\chi^2$ distribution with two degrees of freedom. We convert each $z$ to its corresponding cumulative probability under this distribution and plot the empirical frequency with which these probabilities fall below given thresholds (y-axis) against the theoretical probabilities (x-axis). Perfect calibration corresponds to the diagonal $x = y$; curves above indicate underconfidence (i.e., overestimated uncertainties), while curves below indicate overconfidence.

Fig. S1 presents calibration curves for our main model variants, with calibration quality quantified by the mean squared error (MSE) from the diagonal. As shown, direct decoding models achieve superior calibration as uncertainty is their sole mechanism for handling low signal-to-noise periods, such as during eye fixations or blinks. In contrast, gated decoding models can also rely on gating for such periods, introducing slight redundancy that degrades uncertainty calibration. However, gated models achieve substantially lower L2 errors (see Table 1, main paper), motivating their adoption for deployment.

All models exhibit mild underconfidence, with predicted uncertainty intervals being wider than the observed error distribution warrants. We consider this conservative behavior preferable to overconfidence: it ensures that stated confidence intervals reliably contain the true pupil location, albeit at the cost of slightly overestimating the interval width. This provides trustworthy uncertainty bounds for potential downstream applications, where underestimating positional uncertainty could compromise system reliability.
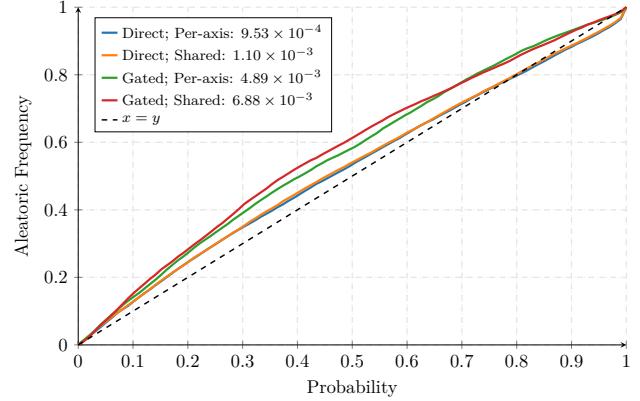


Figure S1. Reliability diagrams for the main model variants. Empirical cumulative frequencies are plotted against theoretical probabilities under a bivariate Gaussian assumption. MSE values quantify calibration quality relative to the diagonal (perfect calibration).
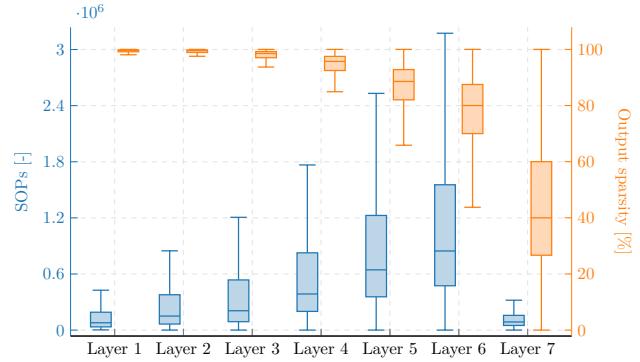


Figure S2. SOPs and output sparsity per layer for the SNN model with gated decoding and shared uncertainty estimation, computed over $\Delta t = 10$ ms windows in simulation.

## S2. Computational load and activation sparsity

Fig. S2 characterizes the computational load and activity patterns across spiking layers on our real-world validation dataset in simulation. Results were obtained with our SNN featuring gated decoding and shared aleatoric uncertainty
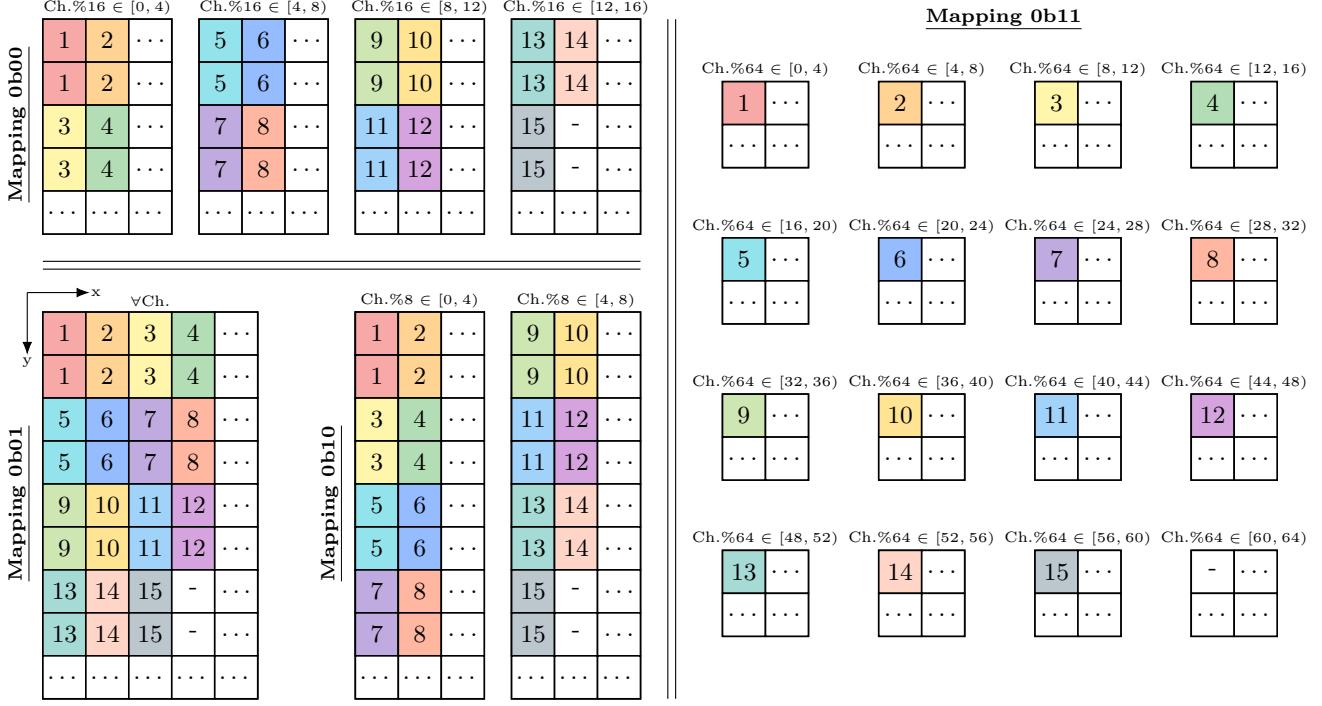
Figure S3. Readout core address mapping characterization. Numbers indicate readout neuron assignments for each feature map location. Repeated numbers denote address collisions; dashes mark inaccessible positions. Ellipsis notation indicates repetition in the corresponding direction: horizontally if aligned right, vertically if aligned below, or diagonally if positioned diagonally.

predictions. The left y-axis (blue) shows SOPs, while the right y-axis (orange) presents layer sparsity, defined as the percentage of spiking neurons remaining silent, both computed using $\Delta t = 10$ ms time windows.

The SOPs distribution reveals that processing demands scale with layer depth, increasing from a median of approx. 0.1M in Layer 1 to 0.8M in Layer 6 before dropping sharply in the final layer. This trend reflects the network's hierarchical feature processing: early layers operate on high-resolution, sparse event inputs, while deeper layers process increasingly abstract representations with denser activation patterns. The dramatic reduction in Layer 7 stems from its role as the output layer with only 15 neurons, minimizing both computational cost and readout bandwidth. All layers operate well below their respective bandwidth thresholds (i.e., 100M SOPs for Layer 1, 30M for Layers 2–7 [53]), validating our activity regularization strategy from Section 3.4 of the main paper.

Sparsity patterns show complementary behavior. Layers 1–3 maintain extremely high sparsity, with most neurons remaining inactive due to the sparse nature of event-based inputs and the feature map resolution. Sparsity gradually decreases through Layers 4–6 as features become more abstract and activation patterns more distributed. Layer 7 exhibits the highest variability (median approx. 40%, range 0–100%), reflecting its role in encoding both pupil location

and aleatoric uncertainty: high sparsity during periods of minimal eye movement or occlusions, and dense activity during active tracking. This adaptive behavior demonstrates that our gated decoding mechanism effectively modulates output activity based on input dynamics, contributing to the event-driven power efficiency shown in Figs. 5 and 6 of the main paper.

## S3. Speck2f's hardware issue

Speck2f's readout core contains a hardware errata [47] affecting neuron addressing when mapping feature maps to its 16 monitoring neurons. Of these, the 0th neuron remains always silent [47], leaving 15 accessible neurons for readout. This core nominally supports multiple spatial configurations, in [channel, height, width] format:

- Mapping 0b00 – Designed for [4, 2, 2] maps.
- Mapping 0b01 – Designed for [2, 2, 4] maps.
- Mapping 0b10 – Designed for [1, 4, 4] maps.
- Mapping 0b11 – Designed for [16, 1, 1] maps.

However, the addressing logic incorrectly resolves spatial coordinates in all configurations. To characterize this hardware errata, we conducted a systematic probing experiment where we stimulated individual neurons across spatial positions and channel indices while monitoring which read-
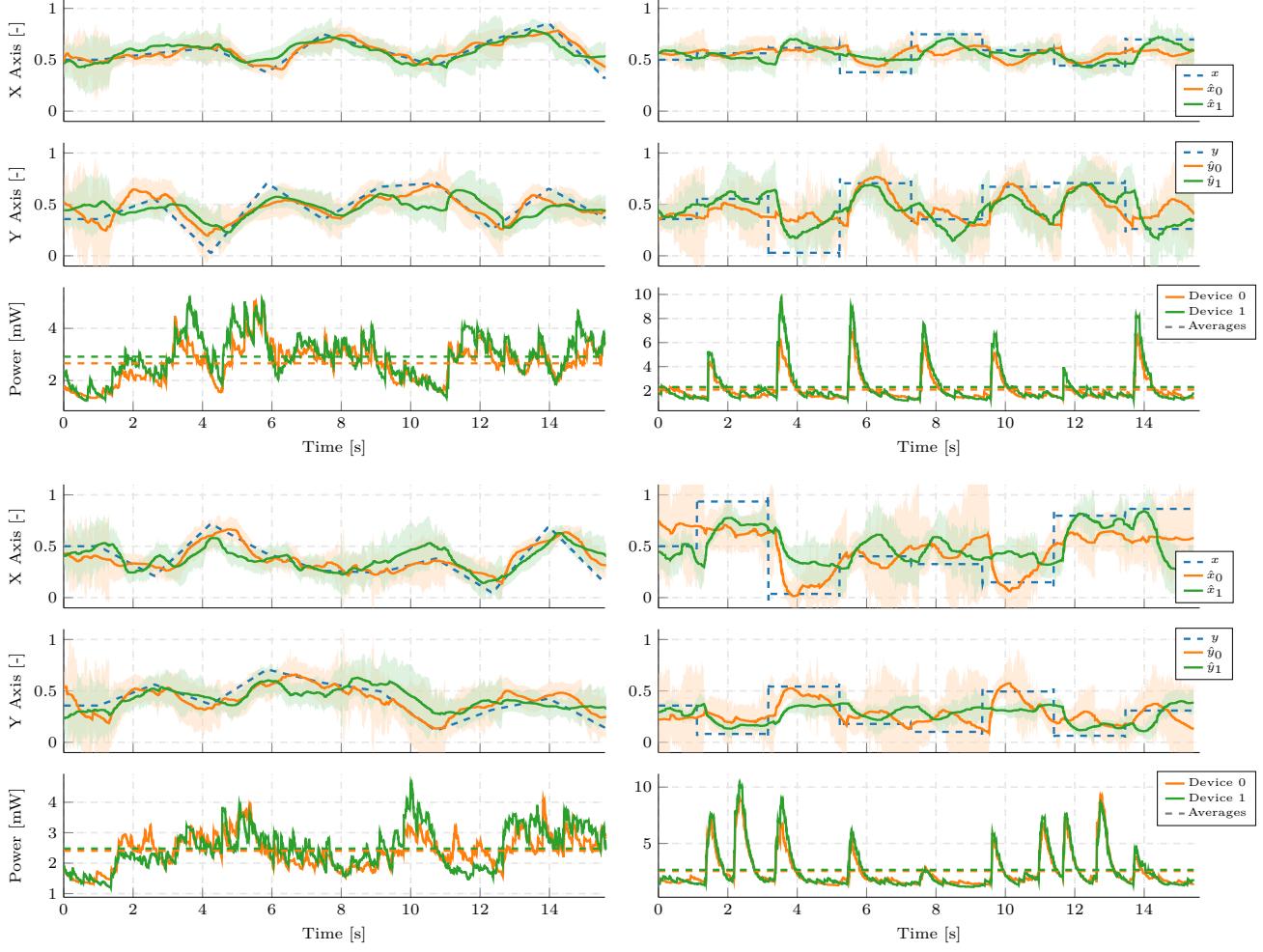
Figure S4. Additional wearable prototype results as a user performs smooth pursuit (left) and saccadic eye movements (right) while tracking a moving dot on a screen at a fixed distance. *Top and middle rows*: Ground-truth and predicted pupil-center coordinates with uncertainty estimates for left (device 0) and right (1) eyes. *Bottom row*: Instantaneous power consumption per device and sequence averages.

out neurons received spikes. Fig. S3 illustrates the resulting address mapping behavior for each configuration. Our characterization reveals systematic address duplication with configuration-dependent patterns:

- Mapping `0b00`: Addresses repeat vertically in pairs and every four positions, horizontally every two positions, and across four-channel groups, requiring a minimal $[13, 4, 2]$ feature map to activate the 15 readout neurons.

- Mapping `0b01`: Addresses repeat vertically in pairs and every eight positions, horizontally every four positions, and across all channels, requiring a minimal $[1, 8, 4]$ feature map to activate the 15 readout neurons.

- Mapping `0b10`: Addresses repeat vertically in pairs and every eight positions, horizontally every two positions, and across four-channel groups, requiring a

minimal $[5, 8, 2]$ feature map to activate the 15 readout neurons.

- Mapping `0b11`: Addresses repeat across four-channel groups only, requiring a minimal $[57, 1, 1]$ feature map to activate the 15 readout neurons.

Our core SNN architecture must route 15 output neurons to the readout core while amplifying their spike counts by a factor of 16 to compensate for the quantized moving average (see Section 3.5, main paper). Based on the results of our hardware characterization, we adopt mapping configuration `0b11` preceded by a $[60, 1, 1]$ layer. This configuration exploits the four-channel grouping pattern inherent to the errata: each group of four consecutive channels (e.g., channels 0–3, 4–7, etc.) maps to the same readout neuron, effectively providing a 4× amplification. This, appended to our proposed spike multiplier (see Section 3.5, main paper)
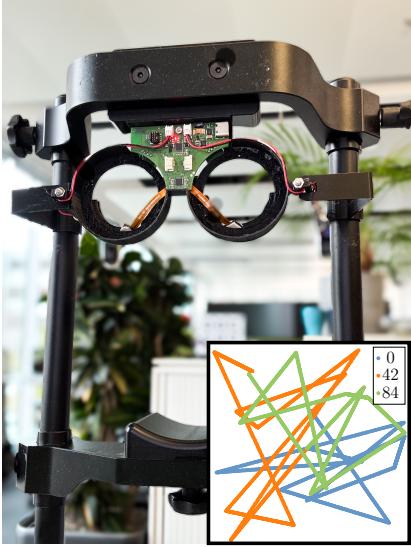
14

Figure S5. Our wearable prototype mounted on a chinrest for controlled data collection. *Inset* shows stimulus trajectories used for both smooth pursuit and saccadic eye movement experiments. Legend indicates random seed values for trajectory generation.

which already contributes a 4× amplification, achieves the required 16× increase in output activity.

## S4. Additional results with wearable prototype

We provide additional wearable prototype results to complement the main paper. Fig. S4 shows representative tracking sequences under various conditions, reinforcing the conclusions from Section 4.3 of the main paper. Across all sequences, we observe consistent performance during smooth pursuit with tight uncertainty bounds and accurate trajectory following, while power consumption remains stable at approx. 2.5 mW per Speck2f device. During saccadic movements, we consistently observe larger overall errors with uncertainty peaking during fixations due to the sparse input events, alongside characteristic power spikes of approx. 5–10 mW corresponding to event bursts. The physical experimental setup and stimulus trajectory patterns used for these experiments are shown in Fig. S5.

### S4.1. Supplementary video

Alongside this document, we provide a video[7] demonstrating real-time operation of our wearable prototype. The video shows live tracking performance, uncertainty estimation, and dynamic power scaling during natural eye movements. We encourage readers to watch it for a complete understanding of our system's capabilities. A snapshot of the real-time monitoring interface is shown in Fig. S6.

---

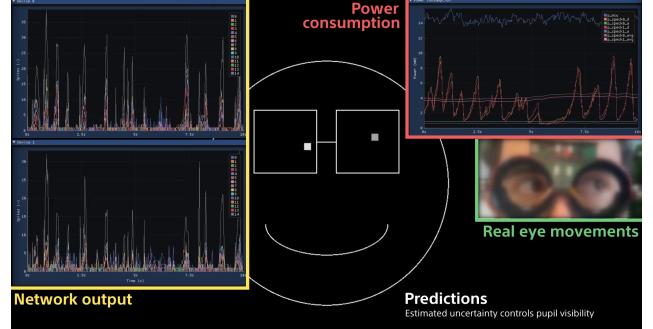[7]The video is available at https://bit.ly/4io51iL



Figure S6. Real-time monitoring interface for the wearable prototype. The interface displays predicted pupil center locations with uncertainty estimates (*center*), power consumption of individual Speck2f devices and MCU (*top right*), and raw SNN output activity (*left*). Data is streamed via Bluetooth at approx. 50 Hz. *Inset* shows corresponding eye movements captured by an external camera for validation (blurred for anonymity).

## S5. Speck2f operation

This section describes the communication protocols and initialization procedures for interfacing with Speck2f without the its FPGA-based development kit. These details complement the official Speck2f documentation [47, 48] to enable replication of our setup.

### S5.1. SPI communication

For each register access, a three-byte command header is constructed prior to SPI transmission to specify the operation type and target register address. The first byte encodes the access mode, where the most significant bit is set to 1 for write operations and cleared for read operations. The following two bytes represent the 16-bit target register address, transmitted in big-endian order. After this header, any associated payload bytes (i.e., data to be written or space reserved for incoming data in a read operation) are transmitted sequentially as part of the same SPI frame.

### S5.2. Reformatting the configuration file

This section presents a systematic approach for compressing and reformatting the binary configuration file generated by SynSense's Samna API [47] for programming the Speck2f via custom SPI-based communication. Tables S1 and S2 summarize important constants and addressing information.

**Kernel memory size computation:** For each layer *l*, we first extract configuration parameters from the original binary file (denoted "binary"). For the kernel size:

$$k_s(l) = (\text{binary}[\text{CNN\_REG}[l] + 0\text{x}1] \gg 2) \,\&\, 0\text{xF} \quad (15)$$

where $\gg$ and $\&$ denote bitwise right-shift and AND.

| Constant | Value |
|---|---|
| KERNEL_START | 0x4C0 |
| R_START | 0x100 |
| KERNEL | [0x10000, 0x14000, 0x18000, |
| | 0x20000, 0x28000, 0x30000, |
| | 0x40000, 0x50000, 0x54000] |
| CNN_REG | [0x280, 0x2C0, 0x300, 0x340 |
| | 0x380, 0x3C0, 0x400, 0x440, |
| | 0x480] |

Table S1. Constants and addresses required during the preparation phase of the Speck2f configuration file.

| Reg. addr. | Config. addr. | Num. regs. |
|---|---|---|
| 0x0000 | 0x0040 | 34 |
| 0x0600 | 0x0280 | 17 |
| 0x0700 | 0x02C0 | 17 |
| 0x0800 | 0x0300 | 17 |
| 0x0900 | 0x0340 | 17 |
| 0x0A00 | 0x0380 | 17 |
| 0x0B00 | 0x03C0 | 17 |
| 0x0C00 | 0x0400 | 17 |
| 0x0D00 | 0x0440 | 17 |
| 0x0E00 | 0x0480 | 17 |
| 0x0100 | 0x0080 | 17 |
| 0x0118 | 0x0098 | 6 |
| 0x0300 | 0x0100 | 45 |
| 0x0200 | 0x00C0 | 23 |

| Kernel addr. |
|---|
| 0x220000 |
| 0x240000 |
| 0x260000 |
| 0x280000 |
| 0x2A0000 |
| 0x2C0000 |
| 0x2E0000 |
| 0x300000 |
| 0x320000 |

Table S2. Register (left) and memory information (right) used in the configuration encoding process.

The output feature count is reconstructed from its low and high byte components:

$$C_{\text{out\_low}}(l) = \text{binary}[\text{CNN\_REG}[l] + 0x04] \gg 6 \quad (16)$$

$$C_{\text{out\_high}}(l) = \text{binary}[\text{CNN\_REG}[l] + 0x05] \quad (17)$$

$$C_{\text{out}}(l) = (C_{\text{out\_high}}(l) \ll 2) \,|\, C_{\text{out\_low}}(l) \quad (18)$$

where $\ll$ and $|$ denote bitwise left-shift and OR.

Similarly, the input feature count is obtained as:

$$C_{\text{in\_low}}(l) = \text{binary}[\text{CNN\_REG}[l] + 0x02] \quad (19)$$

$$C_{\text{in\_high}}(l) = \text{binary}[\text{CNN\_REG}[l] + 0x03]\&0b11 \quad (20)$$

$$C_{\text{in}}(l) = (C_{\text{in\_high}}(l) \ll 8) \,|\, C_{\text{in\_low}}(l) \quad (21)$$

The kernel memory size is then computed as:

$$k_{\text{mem}}(l) = \left((k_s(l) + 1)^2 - 1\right) \ll (w_{\text{in}}(l) + w_{\text{out}}(l))$$
$$|\,(C_{\text{out}}(l) \ll w_{\text{in}}(l)) \,|\, C_{\text{in}}(l) + 1 \quad (22)$$

where $w_{\text{in}}(l) = \lceil \log_2(C_{\text{in}}(l)+1) \rceil$ and $w_{\text{out}}(l) = \lceil \log_2(C_{\text{out}}(l)+1) \rceil$ are the required bit widths.

**File compression:** We run a compression process that iterates through the nine potentially active layers (i.e., $l = 0$ to 8), extracting and packing configuration data. The compressed file begins with binary[0 : KERNEL_START], preserving layer shapes and spiking thresholds. For each layer, we determine if it is active as follows:

$$\text{dest}(l) = \text{binary}[\text{R\_START} + 0x09 + \lfloor l/4 \rfloor] \quad (23)$$

$$\text{layer\_active\_mask} = \text{binary}[\text{R\_START} + 0x27] \quad (24)$$

$$\text{is\_active}(l) = (\text{dest}(l) \gg (2(l \bmod 4))) \,\&0b11$$
$$\vee \, (\text{layer\_active\_mask} \gg l) \,\&0x01 \quad (25)$$

where $\vee$ denotes logical OR.

For active layers, we append: (1) 1-byte layer ID $l$, (2) 2-byte $k_{\text{mem}}(l)$ in little-endian format, (3) $k_{\text{mem}}(l)$ bytes of kernel weights from binary[KERNEL[$l$]], and (4) 2-byte zero for unused leak data. Inactive layers are represented as [$l$, 0, 0, 0, 0].

**File encoding:** The compressed file is converted into SPI command streams for device programming through register and memory configuration operations.

The register configuration process extracts layer parameters from the compressed configuration file and prepares SPI write commands following the protocol in Section S5.1. For each entry $(a_{\text{reg}}, a_{\text{config}}, n)$ in Table S2 (left), $n$ consecutive bytes are read from address $a_{\text{config}}$ and written to Speck2f registers $a_{\text{reg}}$ through $a_{\text{reg}} + n - 1$. Each register write is prepended with a 4-byte big-endian length prefix, enabling SPI controllers to parse transaction boundaries and sequentially transmit individual commands without delimiters or fixed packet sizes.

The memory configuration process prepares the transfer of kernel weights to the Speck2f's on-chip memory banks. For each active layer in the compressed file (starting at address KERNEL_START), if $k_{\text{mem}}(l) > 0$, the kernel data is prepared for writing to memory address mem_addr[$l$] from Table S2 (right). Each memory write requires a two-phase protocol. First, six registers are configured to specify the target memory address and transfer length:

$$0x6 \leftarrow (\text{mem\_addr}[l] \gg 16)\&0xFF$$
$$0x5 \leftarrow (\text{mem\_addr}[l] \gg 8)\&0xFF$$
$$0x4 \leftarrow \text{mem\_addr}[l]\&0xFF$$
$$0x3 \leftarrow (k_{\text{mem}}(l) \gg 16)\&0xFF$$
$$0x2 \leftarrow (k_{\text{mem}}(l) \gg 8)\&0xFF$$
$$0x1 \leftarrow k_{\text{mem}}(l)\&0xFF$$

where, as before, each register write follows the standard SPI protocol plus a 4-byte big-endian length prefix.

Second, the bulk data transfer is initiated by constructing a memory write transaction with length $3 + k_{\text{mem}}(l)$ (i.e., 3-byte command header plus payload). This length is encoded as a 4-byte big-endian prefix, and is followed by the com-

mand header [0xE0, 0x00, 0x00] that identifies a memory write operation, and finally the $k_{\text{mem}}(l)$ kernel weight bytes. This sequence is repeated for each active layer. The final encoded bytearray contains all register and memory write commands, each with its length prefix, ready for sequential SPI transmission.

**Device initialization:** After programming the configuration, the Speck2f requires initialization via a sequence of control register writes. The initialization sequence writes commands 0x01, 0x81, 0xC1, 0xE1, and 0xF1 to register address 0x0 following the SPI protocol in Section S5.1. These commands can be appended to the encoded configuration bytearray for automatic device startup after programming.

### S5.3. Cyclic readout

To monitor spike activity across all 16 readout neurons, we propose a cyclic readout mechanism that sequentially accesses each neuron's spike counter (see Section 3.5, main paper). For each neuron, we toggle the SCLK GPIO pin to trigger a clock cycle, then perform three SPI transmissions: two read operations to retrieve the 16-bit spike count from registers 0x213 (low byte) and 0x214 (high byte), followed by a write to register 0x20B that configures the next neuron for monitoring. The next neuron index is encoded as ((next_neuron & 0xFF) << 2) | 0x83 to ensure proper alignment and command formatting.